

# **FY2005 Accomplishments**

## **NERSC**

## **Architecture Benchmarking and Evaluation**

Lenny Oliker,\* Julian Borrill, Andrew Canning, John Shalf, Jonathan Carter,  
and David Skinner, Lawrence Berkeley National Laboratory  
Stephane Ethier, Princeton Plasma Physics Laboratory

### **Summary**

*LBNL Computing Sciences staff members are applying their expertise in running scientific codes and evaluating HPC system performance to achieve “real world” assessments of leading supercomputers around the world to determine which architectures are best suited for advancing computational science.*

With the re-emergence of viable vector computing systems such as the Earth Simulator and the Cray X1, and with IBM and DOE’s BlueGene/L taking the number one spot on the TOP500 list of the world’s fastest computers, there is renewed debate about which architecture is best suited for running large-scale scientific applications.

In order to cut through conflicting claims, a team of researchers from Berkeley Lab’s Computational Research (CRD) and NERSC Center divisions have been putting various architectures through their paces, running benchmarks as well as scientific applications key to DOE research programs. The team includes Lenny Oliker, Julian Borrill, Andrew Canning, and John Shalf of CRD; Jonathan Carter and David Skinner of NERSC; and Stephane Ethier of Princeton Plasma Physics Laboratory.

As a result of their evaluations, team members have published or had accepted six papers in journals and at conferences in the United States, Norway, Japan, and Spain.

In the initial part of their study, the team traveled to Japan in December 2004, and put

five different systems through their paces, running four different scientific applications key to DOE research programs. As part of the effort, the group became the first international team to conduct a performance evaluation study of the 5,120-processor Earth Simulator. The team also assessed the performance of:

- the 6,080-processor IBM Power3 supercomputer running AIX 5.1 at the NERSC Center
- the 864-processor IBM Power4 supercomputer running AIX 5.2 at Oak Ridge National Laboratory (ORNL)
- the 256-processor SGI Altix 3000 system running 64-bit Linux at ORNL
- the 512-processor Cray X1 supercomputer running UNICOS at ORNL.

“This effort relates to the fact that the gap between peak and actual performance for scientific codes keeps growing,” said team leader Lenny Oliker. “Because of the increasing cost and complexity of HPC systems, it is critical to determine which classes of applications are best suited for a given architecture.”

---

\* 510-486-6625, LOliker@lbl.gov

The four applications and research areas selected by the team for the evaluation were:

- Cactus, an astrophysics code that evolves Einstein's equations from the Theory of Relativity using the Arnowitt-Deser-Misner method
- GTC, a magnetic fusion application that uses the particle-in-cell approach to solve non-linear gyrophase-averaged Vlasov-Poisson equations
- LBMHD, a plasma physics application that uses the Lattice-Boltzmann method to study magnetohydrodynamics
- PARATEC, a first-principles materials science code that solves the Kohn-Sham equations of density-functional theory to obtain electronic wave functions.

"The four applications successfully ran on the Earth Simulator with high parallel efficiency," Olikar said. "And they ran faster than on any other measured architecture — generally by a large margin." However, Olikar added, only codes that scale well and are suited to the vector architecture may be run on the Earth Simulator. "Vector architectures are extremely powerful for the set of applications that map well to those architectures," Olikar said. "But if even a small part of the code is not vectorized, overall performance degrades rapidly."

"We're at a point where no single architecture is well suited to the full spectrum of scientific applications," Olikar added. "One size does not fit all, so we need a range of systems. It's conceivable that future supercomputers would have heterogeneous architectures within a single system, with different sections of a code running on different components."

One of the codes the group intended to run in this study — MADCAP, the Microwave

Anisotropy Dataset Computational Analysis Package — did not scale well enough to be used on the Earth Simulator. MADCAP, developed by Julian Borrill, is a parallel implementation of cosmic microwave background map-making and power spectrum estimation algorithms. Since MADCAP has high I/O requirements, its performance was hampered by the lack of a fast global file system on the Earth Simulator.

Undeterred, the team re-tuned MADCAP and returned to Japan to try again. The resulting paper found that the Cray X1 had the best runtimes but suffered the lowest parallel efficiency. The Earth Simulator and IBM Power3 demonstrated the best scalability, and the code achieved the highest percentage of peak on the Power3. The paper concluded, "Our results highlight the complex interplay between the problem size, architectural paradigm, interconnect, and vendor-supplied numerical libraries, while isolating the I/O filesystem as the key bottleneck across all the platforms."

As for BlueGene/L, currently the world's fastest supercomputer, David Skinner is serving as Berkeley Lab's representative to a new BlueGene/L consortium led by Argonne National Laboratory. This consortium will work together to develop or port BlueGene applications and system software, conduct detailed performance analysis on applications, develop mutual training and support mechanisms, and contribute to future platform directions.

**For further information on this subject contact:**

David Goodwin, NERSC Program Manager  
Mathematical, Information, and Computational  
Sciences Division  
Office of Advanced Scientific Computing Research  
Phone: 301-903-6474  
dave.goodwin@science.doe.gov

## **NERSC Develops Archiving Strategies to Help Genome Researchers Sort Through Billions of Data Files**

Nancy Meyer,\* Harvard Holmes, Damian Hazen, and Wayne Hurlbert,  
Lawrence Berkeley National Laboratory  
Alex Copeland and Yunian Lou, Joint Genome Institute

### **Summary**

*When researchers at the Production Genome Facility at DOE's Joint Genome Institute found they were generating data faster than they could find room for the files, let alone make them easily accessible for analysis, a collaboration with NERSC's Mass Storage Group developed strategies for improving the reliability of data storage while also making retrieval easier.*

DOE's Joint Genome Institute (JGI) is one of the world's leading facilities in the scientific quest to unravel the genetic data that make up living things. With advances in automatically sequencing genomic information, scientists at the JGI's Production Genome Facility (PGF) found themselves overrun with sequence data as their production capacity had grown so rapidly that data had overflowed the existing storage capacity. Since the resulting data are used by researchers around the world, ensuring the data are both reliably archived and easily retrievable are key issues.

As one of the world's largest public DNA sequencing facilities, the PGF produces 2 million files per month of trace data (25 to 100 KB each), 100 assembled projects per month (50 MB to 250 MB), and several very large assembled projects per year (~50 GB). The total averages ~2000 GB per month.

In addition to the amount of data, a major challenge is that way the data are produced. Data from the sequencing of many different organisms are produced in parallel each day, such that a daily archive spreads the data for

a particular organism over many tapes. Current sequencing methods generate a large volume of trace files that have to be managed — typically 100,000 files or more. And to check for errors in the sequence or make detailed comparisons with other sequences, researchers often need to refer back to these traces. Unfortunately, these traces are usually provided as a group of files with no information on where the traces occur in the sequence.

This problem was compounded by the PGF's lack of sufficient online storage, which made organization (and subsequent retrieval) of the data difficult and led to unnecessary replication of files. This situation required significant staff time to move files and reorganize file systems to find sufficient space for ongoing production needs; and it required auxiliary tape storage that was not particularly reliable.

Staff from NERSC's Mass Storage Group and the PGF agreed to work together to address two key issues. The most immediate goal was to for NERSC High Performance Storage System (HPSS) to become the ar-

---

\* 510-486-6627, NLMeyer@lbl.gov

chive for the JGI data, replacing the less-reliable local tape operation and freeing up disk space at the PGF for more immediate production needs. The second goal was to collaborate with JGI to improve the data handling capabilities of the genome sequencing and data distribution processes.



*NERSC's HPSS storage system is robust and available 24 hours a day, seven days a week, as well as highly scalable and configurable. NERSC also has high-quality, high-bandwidth connectivity to the other DOE laboratories and major universities provided by ESnet.*

Most of the low-level data produced by the PGF are now routinely archived at NERSC, with ~50 GB of raw trace data being transferred from JGI to NERSC each night. This archive of data also forms a foundation for consideration of further steps to enhance the utility of the data.

To accomplish the archive process, NERSC staff came up with the following solutions to address the main challenges:

- the use of an HPSS variant of tar files (HTAR) to combine multiple small files into chunks large enough for efficient transfer and storage
- the design and implementation of a directory structure that would allow easy location of the various files

- the creation of scripts that would run on the PGF machines to transfer the files
- network tuning and configuration changes to support and optimize the data transfer between the PGF and NERSC.

By using these techniques, the archiving system can be scaled up over time as the amount of data continues to increase — up to billions of files can be handled with these techniques. The data have been aggregated into larger collections which hold tens of thousands of files in a single file in the NERSC storage system. This data can now be accessed as one large file, or each individual file can be accessed without retrieving the whole aggregate.

Not only will the new techniques be able to handle future data, they also helped when the PGF staff discovered raw data that had previously been processed by software that had an undetected bug. The staff were able to retrieve the raw data from NERSC and reprocess it in about 1½ months, rather than go back to the sequencing machines and produce the data all over again, which would have taken about six months. In addition to saving time, this also saved money — a rough estimate is that the original data collection comprised up to 100,000 files/day at a cost of \$1 per file, which added up to \$1.2 million for processing six months' worth of data. Comparing this figure to the cost of a month and a half of staff time, the estimated savings are about \$1 million — and the end result is a more reliable archive.

**For further information on this subject contact:**

David Goodwin, NERSC Program Manager  
Mathematical, Information, and Computational  
Sciences Division  
Office of Advanced Scientific Computing Research  
Phone: 301-903-6474  
dave.goodwin@science.doe.gov

## **NERSC Reaches Another Checkpoint/Restart Milestone**

William T.C. Kramer\*, James Craw, and Jay Srinivasan,  
Lawrence Berkeley National Laboratory

### **Summary**

*In 1997 NERSC made history by being the first computing center to achieve successful checkpoint/restart on a massively parallel system, the Cray T3E. In 2005, NERSC and IBM achieved the first full-scale use of checkpoint/restart software with an actual production workload on an IBM SP, as well as the first checkpoint/restart on a system with more than 2,000 processors.*

Checkpointing—successfully stopping and restarting a number of jobs on a supercomputer without any data processing loss or discontinuity—is an important productivity capability for a scientific computing center.

Checkpointing maximizes system availability for users and minimizes wasted compute cycles because no recomputation is necessary after restarting. Checkpointing means stopping a program in progress and saving the current state of the program and its data—in effect, “bookmarking” where the program left off so it can start up later in exactly the same place. The process records all the information, transfers that information out of the machine, then puts information back in and gets it all running again with no loss of processing time or data. Recovery of the unfinished applications resumes from the point of interruption.

Although being able to stop and restart a computer system without data loss is important for any system, the value is much greater as the size of the system increases. For example, without checkpointing, when a

single-processor system runs 12 hours of computing work, is interrupted and cannot be restarted, it loses 12 hours of work. On the other hand, when a 2,000-processor system runs 12 hours and cannot be restarted, it loses 24,000 hours of computing work. Checkpointing on massively parallel systems is quite difficult because of the complexity of synchronizing so many processors.

In 1997 NERSC made history by being the first computing center to achieve successful checkpoint/restart on a massively parallel system, the Cray T3E. The weekend of June 11–12, 2005 marked another milestone when IBM personnel used NERSC’s 6,656-processor IBM supercomputer, Seaborg, for dedicated testing of IBM’s latest HPC Software Stack, a set of tools for high performance computing.

To maximize system utilization for NERSC users, instead of “draining” the system (letting running jobs continue to completion) before starting this dedicated testing, NERSC staff checkpointed all running jobs at the start of the testing period. This is believed to be the first full-scale use of the

---

\* 510-486-7577, WTKramer@lbl.gov



checkpoint/restart software with an actual production workload on an IBM SP, as well as the first checkpoint/restart on a system with more than 2,000 processors.



*Seaborg, NERSC's IBM supercomputer, is even more productive with checkpoint/restart capability, which prevents loss of data during system interruptions.*

This achievement is the culmination of a collaborative effort between NERSC and IBM that began in 1999. Over the past two years, NERSC staff worked with IBM to resolve a significant number of technical and design issues before checkpoint/restart was put into production use. Additionally, NERSC staff made changes to the job submission mechanism used on Seaborg and developed a custom program to enable all jobs to be checkpointable before they are submitted to the LoadLeveler job scheduling system.

Of the 44 jobs that were checkpointed during the test, approximately 65% checkpointed successfully. Of the 15 jobs that did not checkpoint successfully, only 7 jobs were deleted from the queuing system, while the rest were requeued to run again at a later time. This test enabled NERSC and IBM staff to identify some previously undetected problems with the checkpoint/restart software, which they were then able to correct.

The HPC Software Stack changes were transparent to NERSC users, so they did not notice any difference in how the batch system accepted, queued, and ran their jobs. However, most jobs submitted to LoadLeveler are now checkpointable, which allows NERSC staff to use checkpoint/restart to migrate jobs from nodes and shorten queue drain times before system or node maintenance.

Checkpoint/restart for LoadLeveler is now in production use, allowing NERSC staff to minimize disruptions and downtime for more than a thousand users around the country, making NERSC an even more valuable computational science resource.

**For further information on this subject contact:**  
David Goodwin, NERSC Program Manager  
Mathematical, Information, and Computational  
Sciences Division  
Office of Advanced Scientific Computing Research  
Phone: 301-903-6474  
[dave.goodwin@science.doe.gov](mailto:dave.goodwin@science.doe.gov)

## **NERSC Launches Two New Systems to Advance Scientific Computing**

William T.C. Kramer, \* Lawrence Berkeley National Laboratory

### **Summary**

*NERSC has put two new computing systems into production: The Jacquard system is one of the largest production InfiniBand-based Linux cluster systems and has met unprecedented acceptance criteria for performance, reliability and functionality. DaVinci is a visualization and data analysis server that offers interactive access to large amounts of memory and high performance I/O capabilities typically required to analyze large datasets.*

In August 2005, the U.S. Department of Energy's National Energy Research Scientific Computing Center (NERSC) put two new computing systems into production, providing additional computing capabilities for NERSC's 2,500 researchers at national labs and universities across the country.

On Monday, Aug. 1, NERSC announced that the new Linux Networkx cluster system had entered full production mode for NERSC users. The new system has a total of 702 processors, including 640 processors for computing and 40 for storage.

Named "Jacquard," the Linux Networkx system will provide computational resources to scientists who run jobs on up to 124 processors, freeing up more resources on Seaborg for jobs which scale to 512 or more processors.

The Jacquard system is one of the largest production InfiniBand-based Linux cluster systems and has met rigorous acceptance criteria for performance, reliability and functionality that are unprecedented for an InfiniBand-based cluster. Jacquard is the first system to deploy Mellanox 12X InfiniBand uplinks in its fat-tree interconnect, reducing network hot spots and improving reliability

by dramatically reducing the number of cables required.



*Jacquard will run jobs using up to 124 processors, freeing up resources on Seaborg for larger-scale jobs.*

The system has 640 AMD 2.2 GHz Dual Opteron™ processors devoted to computation, with the rest used for I/O, interactive work, testing and interconnect management. Jacquard has a peak performance of 3.1 trillion floating point operations per second (teraflop/s). Storage from DataDirect Networks provides 30 terabytes of globally available formatted storage.

The acceptance test included a 14-day availability test during which a select group of

\* 510-486-7577, WTKramer@lbl.gov



NERSC users were given full access to the Jacquard cluster to thoroughly test the entire system in production operation. Jacquard had a 99 percent availability uptime during the testing while users and scientists ran a variety of codes and jobs on the system. The thorough acceptance testing by NERSC ensures Jacquard is ready for a production environment for thousands of scientists and researchers across the nation.

The installation also includes a smaller development cluster called Jacdev. This system consists of 20 processors.

Following the tradition at NERSC, the system was named for someone who has had an impact on science and/or computing. In 1801, Joseph-Marie Jacquard invented the Jacquard loom, which was the first programmable machine. The Jacquard loom used punched cards and a control unit that allowed a skilled user to program detailed patterns on the loom.

In mid-August, NERSC put into production a new server specifically tailored to interactive visualization and data analysis work. The 32-processor SGI Altix, called DaVinci, offers interactive access to large amounts of memory and high performance I/O capabilities typically required to analyze the large datasets produced by the NERSC high performance computing systems (Jacquard and Seaborg).

“With its 192 gigabytes of RAM and 25 terabytes of disk, DaVinci’s balance is biased toward memory and I/O, which is different from the other systems at NERSC,” said John Shalf of LBNL’s Visualization Group. “This design gives us expanded capabilities for data analysis and analytics, along with interactive visualization.”

DaVinci, named for the Renaissance artist, scientist, writer and engineer, has 6 gigabytes of memory per processor, compared to 1 gigabyte per processor on Seaborg and 4 gigabytes on the new Linux Networx cluster, Jacquard.

Users can get interactive access to all 192 gigabytes of memory from a single application, whereas the interactive limits on production NERSC supercomputing systems restrict interactive tasks to a far smaller amount of memory.

The new server will run a number of visualization, statistics and mathematics applications including IDL, Mathematica, Star-P (a parallel implementation of MatLab), AVS/Express, LLNL VisIT (a parallel visualization application), and CEI Ensignt. Many users depend on IDL and MatLab to process or reorganize data in preparation for visualization. The large memory will particularly benefit these types of jobs.

While DaVinci will be available for interactive use by day, by night the system will be set up to run batch jobs, especially those jobs that are data intensive.

DaVinci is already connected to the HPSS and ESnet networks at NERSC by two independent 10 gigabit Ethernet connections and is expected to be integrated in the Facility Wide File Sharing System later this year.

**For further information on this subject contact:**  
David Goodwin, NERSC Program Manager  
Mathematical, Information, and Computational  
Sciences Division  
Office of Advanced Scientific Computing Research  
Phone: 301-903-6474  
dave.goodwin@science.doe.gov